

# Bayesian Model Averaging

SIBGHAT & SIDDHANT

# Background Problem

- ▶ Model Uncertainty
- ▶ Model Selection
- ▶ Better Prediction

# BMA Introduction: Notations

- ▶  $\Delta$  is quantity of interest
- ▶ Future Observable
- ▶ Utility of a course of action
- ▶  $D$  is data
- ▶  $\mathbf{M} = \{M_k, k=1,2,3,\dots,K\}$
- ▶  $\Theta_k$  vector of parameters in  $M_k$
- ▶  $\Pr(\Theta_k \mid M_k)$  prior density of  $\Theta_k$  under  $M_k$
- ▶  $\Pr(D \mid \Theta_k, M_k)$  is likelihood for data
- ▶  $\Pr(M_k)$  is prior probability that  $M_k$  is the true Model

# Mathematics for BMA

- Posterior distribution given data  $D$  is

$$\Pr(\Delta \mid D) = \sum_{k=1}^K \Pr(\Delta \mid M_k, D) \Pr(M_k \mid D)$$

Posterior probability for model  $M_k \in \mathbf{M}$  is

$$\Pr(M_k \mid D) = \frac{\Pr(D \mid M_k) \Pr(M_k)}{\sum_{l=1}^K \Pr(D \mid M_l) \Pr(M_l)}$$

where

$$\Pr(D \mid M_k) = \int \Pr(D \mid \theta_k, M_k) \Pr(\theta_k \mid M_k) d\theta_k$$

# Expectation using BMA

► Let  $\Delta^k = E[\Delta \mid D, M_k]$

► Posterior mean of  $\Delta$  using BMA is

$$E[\Delta \mid D] = \sum_{k=1}^K \Delta^k \text{pr}(M_k \mid D)$$

# Implementation Issues

- ▶ Simple right ?
- ▶ Not so Much....

# Issues ...

- ▶ BMA is nice but difficult to implement
- ▶ Reasons:
  - $M$  can be enormous. Infeasible to sum over all models
  - Integrals can be hard to compute, even using MCMC
  - What about Prior distribution ?
  - Is it worth it ?

# Managing the Summation

- Feasible way to compute the equation

$$\sum_{k=1}^K \Pr(\Delta \mid Mk, D) \Pr(Mk \mid D)$$

Two approaches

- Occam's window method
- Monte Carlo Markov Chains Model Composition



# Occam's window method

- ▶ Average over a subset of models supported by the data
- ▶ Principle 1
- ▶ Disregard a model if it predicts the data far less than the model with best predictions
- ▶ Formally:

$$A' = \{M_k; \frac{\max_l (\Pr(M_l | D))}{\Pr(M_k | D)} \leq C\}$$

# Occam's Window method ...

- ▶ Exclude complex model if the data supports the simpler model (Occam's razor)
- ▶ Formally

$$B = \{Mk : \exists Ml \in A', Ml \subset Mk, \Pr(Ml | D) / \Pr(Mk | D) > 1\}$$

Subset of models to be used is :  $A = A' \setminus B$

All probabilities conditional on A

# Use MCMC Model composition

- ▶ Use MCMC to directly approximate the first equation
- ▶ Construct a Markov Chain to  $\{M(t)\}$ ,  $t=1,2,\dots$  with state space  $\mathbf{M}$  and stationary distribution  $\Pr(M_i \mid D)$
- ▶ Simulate chain to get observations  $M(1), \dots, M(N)$
- ▶ Then for any function  $g(M_i)$  defined on  $\mathbf{M}$ , compute average

- ▶ 
$$\mathbf{G}(\text{est}) = \frac{\sum_{t=1}^N g(M(t))}{N}$$

# Computing Integrals

- Integrals of the form:

$$\mathbf{Pr}(\mathbf{D} | \mathbf{Mk}) = \int \mathbf{Pr}(\mathbf{D} | \theta_k, \mathbf{Mk}) \mathbf{pr}(\theta_k | \mathbf{Mk}) d\theta_k$$

Can be difficult to compute.

Solution ?

# Computing Integrals...

- ▶ Closed form integrals available for multiple Regression & graphical models
- ▶ Laplace method helps approximate  $\Pr(D \mid M_k)$  & sometimes yields BIC approximation
- ▶ Approximate  $\Pr(\Delta \mid M_k, D)$  with  $\Pr(\Delta \mid M_k, \theta(\text{estimated}), D)$  where  $\theta(\text{estimated})$  is MLE

# BMA for Linear Regression: Predictors, Outliers & Transformations

- Suppose a dependent variable  $Y$  and predicts  $X_1, \dots, X_k$ . Then variable selection methods try to find the “best” model with the form

$$Y = \sum_{j=1}^p B_{ij} X_{ij} + \varepsilon + B_0$$

BMA however tries to average over all possible sets of predictors

# Linear Regression: Transformation & Outliers

- ▶ We can use Box-cox transformation for the response

$$y^{(p)} = \begin{cases} \frac{y^p - 1}{p} & \text{if } p \text{ not equal to zero} \\ \log(y) & \text{if } p \text{ equal to zero} \end{cases}$$

And the model is  $y^{(p)} = XB + \varepsilon$  where  $\varepsilon \sim N(0, \sigma^2 I)$

We can use “change point transformations” to transform the predictors

- Use the output from the alternating conditioning expectation algorithm to suggest the form of transformation
- Use Bayes factor to choose the precise transformation

# Linear Regression: Transformation & Outliers

- ▶ We can use variance-inflation model for outliers by assuming:

$$\varepsilon = \begin{cases} N(0, \sigma^2) & \text{w.p. } (1 - \pi) \\ N(0, K^2 \sigma^2) & \text{w.p. } \pi \end{cases}$$

Simultaneous variable & Outlier selection method

- Use a highly robust technique to identify potential outliers
- Compute all possible posterior model probabilities or use MC3, considering all possible subsets of potential outliers.



# Generalized Linear Models

- The Bayes factor for model  $M_1$  against  $M_0$ :

$$B_{10} = \frac{\Pr(D | M_1)}{\Pr(D | M_0)}$$

Consider  $M+1$  models  $M_0, M_1, \dots, M_K$ . Then the posterior probability for Model  $M_i$  using Bayes Factor is:

$$\Pr(M_i | D) = \frac{\alpha_i B_{i0}}{\sum_{j=0}^K \alpha_j B_{j0}}$$

Where  $\alpha_i = \frac{\Pr(M_i)}{\Pr(M_0)}$   $i = 0, \dots, K$

# Generalized Linear Model

- ▶ Dependent variable:  $Y_i$
- ▶ Independent variables:  $X_i = (x_{i1}, \dots, x_{ip})$ ,  $i=0, \dots, n$  where  $x_{i1}=1$
- ▶ The null model  $M_0$  is defined as  $B_j=0$  where  $j=2, \dots, p$

# Generalized Linear Models ....

- If we use Laplace approximation, we get

$$\Pr(D | Mk) \approx (2\pi)^{p_k/2} |\Psi|^{1/2} \Pr(D | B_k^{\wedge}, Mk) \Pr(B_k^{\wedge} | Mk)$$

Where  $p_k$  is the dimension of  $B_k$ , and  $B_k^{\wedge}$  is the posterior mode of  $B_k$ , and  $\Psi_k$  is minus the inverse hessian of

$$h(B_k) = \{\Pr(D | B_k, Mk) \Pr(B_k | Mk)\} \text{ evaluated at } B_k = B_k^{\wedge}$$

# Generalized Linear Models ...

- ▶ Suppose  $E[B_k | M_k] = w_k$  and  $\text{var}(B_k | M_k) = W_k$
- ▶ Use one step of Newton's method to approximate  $B_k^{\hat{}}$  starting from  $B^{\hat{}}$
- ▶ Then we have the approximation

$$2 \log B_{10} \approx X^2 + (E1 - E0)$$

And

$$X^2 = \{l_1(B_1^{\hat{}}) - l_0(B_0^{\hat{}})\}$$

$$l_k(B_k) = \log\{\text{Pr}(D | B_k, M_k)\}$$

$$E_k = 2\lambda_k(B_k^{\hat{}}) + \lambda'_k(B_k^{\hat{}})^T (F_k + G_k)^{-1} \{2 - F_k (F_k + G_k)^{-1}\} \lambda'_k(B_k^{\hat{}}) - \log |F_k + G_k| + p_k \log(2\pi)$$

Where  $F_k$  is the expected fisher information matrix,  $G_k = (W_k)^{-1}$   
and  $\lambda_k = \log \text{Pr}(B_k | M_k)$

# Survival Analysis

- ▶ Hazard rate  $\lambda(t) = f(t) / 1 - F(t)$
- ▶ Cox proportional hazard model:  $\lambda(t | X_i) = \lambda_0(t) \exp(X_i B)$  where  $\lambda_0(t)$  is the baseline hazard rate at time  $t$ .
- ▶ The estimation of  $B$  is based on partial likelihood.

$$PL(B) = \prod_{i=1}^n \left( \frac{\exp(X_i B)}{\sum_{l \in R_i} \exp(X_l^T B)} \right)^{w_i}$$

Where  $R_i$  is the risk at time  $t_i$  and  $w_i$  indicates whether subject  $i$  is censored or not

# Survival Analysis ...

- ▶ A lot of studies have used the MLE approximation for survival analysis models

$$\Pr(\Delta \mid Mk, D) \approx \Pr(\Delta \mid Mk, \tilde{B}_k, D)$$

And the Laplace approximation

$$\log \Pr(D \mid M) \approx \log \Pr(D \mid Mk, \tilde{B}_k) - d_k \log(n)$$

Where  $d_k$  is the dimension of  $\tilde{B}_k$

# Graphical Models: Missing Data & Auxiliary Variables

- ▶ A graphical model is a statistical model with a set of conditional independence relationships being described by means of a graph.
- ▶ We'd study here only acyclic direct graph
- ▶ Use either analytical or numerical approximations when we apply BMA and Bayesian graphical models to solve problems with missing data. For example

$$\frac{\Pr(D | M0)}{\Pr(D | M1)} = E \left( \frac{\Pr(D, Z | M0)}{\Pr(D, Z | M1)} \mid D, M1 \right)$$

Where Z denotes the missing data and/or auxiliary variables

# Specifying prior model probabilities

- ▶ It has been proved that in BMA informative prior have better predictive performance than the neutral priors
- ▶ Consider the equation

$$\Pr(M_i) = \pi_j^{\delta_{ij}} (1 - \pi_j)^{1 - \delta_{ij}}$$

Where  $M_i$  is a linear model with  $p$  covariates, and  $\pi_j$  is the prior probability that  $B_j \neq 0$  and  $\delta_{ij}$  indicates whether  $X_j$  is included in  $M_i$  or not

- In the case of graphical models, prior probability for each link is specified and multiplied eventually.



# Success of a model/Predictive performance

- ▶ Split data in two halves randomly, called training & test data sets

- ▶ Predictive log score:

Single Model :  $\sum_{d \in \text{test\_data}} -\log \Pr(d | M, \text{training\_data})$

BMA:  $\sum_{d \in \text{test\_data}} -\log \sum_{M \in A} \{ \Pr(d | M, \text{training\_data}) \Pr(M | \text{training\_data}) \}$

Smaller P.L.S indicates better predictive performance

# Philosophical debate

- ▶ Is it Worth it ?
- ▶ Example in R.

# Questions ?

- ▶ Questions regarding BMA